

A Two-Level Direct Solver for the Hierarchical Poincaré-Steklov Method

Joseph Kump^{a,*}, Anna Yesypenko^a, Per-Gunnar Martinsson^a

^a*Oden Institute, University of Texas at Austin, 201 E 24th St, Austin, 78712, TX, USA*

Abstract

We introduce a two-level direct solver for the Hierarchical Poincaré-Steklov (HPS) method for solving linear elliptic PDEs. HPS combines multidomain spectral collocation with a direct solver, enabling high-order discretizations for highly oscillatory solutions while preserving computational efficiency. Our method employs batched linear algebra routines with GPU acceleration to reduce the problem to subdomain interfaces, yielding a block-sparse linear system. This system is then factorized using a sparse direct solver that employs pivoting to achieve better numerical stability than the original HPS scheme. For a discretization of local order p involving a total of N degrees of freedom, the initial reduction step has asymptotic complexity $\mathcal{O}(Np^6)$ in three dimensions. Nevertheless, the high efficiency of batched GPU routines makes the overall cost for practical purposes independent of polynomial order (for order $p = 20$ or even higher). Additionally, the cost of the sparse direct solver is independent of the polynomial order. We present a description and justification of our method, along with numerical experiments on three-dimensional problems to evaluate its accuracy and performance.

Keywords: sparse direct solvers, elliptic PDEs, spectral methods, domain decomposition, GPUs

*Corresponding Author

Email address: josek97@utexas.edu (Joseph Kump)

1. Introduction

The manuscript describes a numerical method for accurately and efficiently solving elliptic boundary value problems of the form

$$\begin{aligned} Au(\mathbf{x}) &= f(\mathbf{x}), & \mathbf{x} \in \Omega \\ u(\mathbf{x}) &= g(\mathbf{x}), & \mathbf{x} \in \partial\Omega, \end{aligned} \quad (1)$$

where Ω is a regular bounded domain in \mathbb{R}^d for $d = 2$ or 3 . Further, A is an elliptic partial differential operator, $g : \partial\Omega \rightarrow \mathbb{R}$ is a given Dirichlet boundary condition, and $f : \Omega \rightarrow \mathbb{R}$ is a given body load. As a concrete example, we focus on the 3D variable-coefficient Helmholtz operator

$$Au = - \left(\frac{\partial^2}{\partial x_1^2} + \frac{\partial^2}{\partial x_2^2} + \frac{\partial^2}{\partial x_3^2} \right) u - \kappa^2 b(\mathbf{x})u. \quad (2)$$

Then (1) models time-harmonic wave phenomena, with κ denoting the wave-number and b a spatially varying, nonnegative function that encodes material heterogeneity.

We introduce a modified solver for the HPS scheme on 2D and 3D domains that forgoes the hierarchical approach of constructing solution operators for subdomain boundaries, instead using a two-level framework. This builds upon work in [1], which was limited to problems in two dimensions. Rather than recursively constructing a chain of solution operators in an upward pass (as in the original HPS method [2]), we use the DtN maps of each lowest-level subdomain ("leaf box") to assemble a block-sparse system that encodes all leaf box boundaries concurrently. By creating a single larger system, we can interface with multifrontal sparse solvers such as MUMPS, which apply pivoting to improve numerical stability [3]. Moreover, the cost of factorizing the sparse system is independent of the choice of subdomain polynomial order p . The DtN maps for the subdomains are constructed in parallel with GPU acceleration, where the use of hardware acceleration mitigates the impact of the local polynomial order p .

GPU acceleration of PDE solvers has received significant attention in recent years, mostly for low-order numerical methods such as in [4, 5, 6, 7]. A recurring theme in these works is using GPUs to reduce the problem along each discretization element to a smaller or more structured linear system. We adopt a similar approach, reducing the problem from one on all HPS discretization points to just subdomain boundary points. This yields a significantly smaller sparse system (acting on $O(p^{d-1})$ discretization points per

subdomain instead of $O(p^d)$) which is easier to store and faster to factorize, and compatible with black-box direct solvers. The uniformity of this reduction – applied identically across subdomains – also enables the use of efficient batched linear algebra routines, allowing us to recompute the required local operators on the fly rather than store them. We assemble the sparse system’s data (the DtN maps) and perform localized leaf box solves with batched linear algebra routines on GPUs via PyTorch [8].

Previous work has sought different means to accelerate either the box interior or box boundary solves present in HPS. For instance, [9] leverages rank structure in 3D HPS operators to compress matrices and apply fast structured matrix algebra, while [10] uses matrix compression to accelerate nested dissection for 2D problems where solutions near the boundary are required. In terms of parallelism, [11] employs OpenMP and MKL to parallelize the upward pass and operator construction within a shared memory framework. More recently, [12] proposed a distributed memory approach, assembling a global sparse system from all discretization points and solving it iteratively using a GMRES iterative solver with MPI parallelism. In contrast, our method constructs a sparse system using only the box boundary degrees of freedom and retains the use of a direct solver.

This paper provides a summary of our two-level HPS solver. Section 2 introduces the block-sparse linear system associated with HPS that we wish to solve, Section 3 describes the complete algorithm for constructing and solving the reduced system, and Section 4 analyzes performance and demonstrates practical applications.

2. Method Overview

In this section we introduce the core components of the method we use to solve the boundary value problem (1). We begin by constructing a sparse linear system based on spectral differentiation within subdomains. This leads to a block-sparse global system, which we simplify by eliminating interior degrees of freedom to expose a smaller system that involves only degrees of freedom associated with discretization nodes on subdomain boundaries.

2.1. Global Discretization

Equation (1) can be solved using the Hierarchical Poincaré-Steklov (HPS) method [2, 13, 14], which combines a multidomain spectral collocation discretization with a nested dissection based direct solver. The discretization

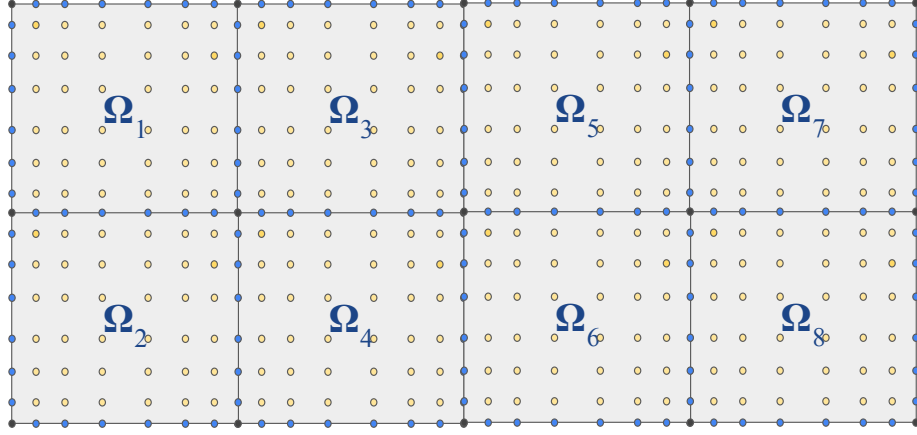


Figure 1: Diagram of a 2D domain partitioned into subdomains for HPS. Yellow nodes are subdomain interiors, blue nodes are subdomain boundaries, and gray nodes are corners.

starts with a flat tessellation into subdomains (“boxes”), where the PDE is enforced locally through spectral collocation in each subdomain, and then continuity of the normal derivative is enforced across boundaries. Like other multidomain spectral methods, HPS enables use of high-order local discretizations that mitigate the pollution effect [15, 16, 17]. It also requires only minimal subdomain overlap, unlike high-order finite difference methods, which reduces computational cost in direct solvers [18, Chapter 24] and inter-process communication in parallel settings [19]. In previous HPS solvers such as [2, 9, 11, 20], the resulting discretized system is then solved via the hierarchical merging of Dirichlet-to-Neumann (DtN) maps.

The domain $\Omega \subset \mathbb{R}^d$ is partitioned into a series of non-overlapping subdomains, often called “leaf boxes” or just “boxes”. Generally these subdomains are rectangular and of the same size, but their shape may be modified through the use of parameter maps (as seen in Section 4.3). An example of a 2D discretization with uniform rectangular subdomains is shown in Figure 1. On each subdomain $\Omega^{(j)}$ we place a tensor product of p^d Chebyshev nodes, with the choice of polynomial order p consistent across all subdomains. These nodes are shared along the boundaries between subdomains: edges in 2D and faces in 3D. We aim to solve for the solution u of our boundary value problem on these discretization points, which we denote as the vector \mathbf{u} .

On the subdomain interiors, we locally enforce the differential equation (1) using a spectral collocation method [21]. On the subdomain boundaries,

we enforce continuity of the Neumann derivative of the two adjacent subdomains. More detailed formulations for how we encode these conditions can be found in [18, Chapter 24]. Together, the matrices for interior and boundary solves form the block rows of a larger system \mathbf{A} that concurrently solves for all subdomain interior and boundary points of our discretization,

$$\mathbf{A}\mathbf{u} = \mathbf{f}. \quad (3)$$

For a discretization node $\mathbf{x}_i \in \Omega$, the value of $[\mathbf{A}\mathbf{u}](i)$ is

$$[\mathbf{A}\mathbf{u}](i) \approx \begin{cases} Au(\mathbf{x}_i) & \text{if } \mathbf{x}_i \in \text{interior of some box } \Omega^{(j)}, \\ \frac{\partial u}{\partial n}(\mathbf{x}_i)|_{\Omega^{(j)}} + \frac{\partial u}{\partial n}(\mathbf{x}_i)|_{\Omega^{(k)}} & \text{if } \mathbf{x}_i \in \text{boundary } \partial\Omega^{(j)} \cap \partial\Omega^{(k)}. \end{cases} \quad (4)$$

Nodes that lie on corners (and edges in 3D) in principle require special treatment, but in most applications they can surprisingly be dropped from consideration entirely. See Remark 1.

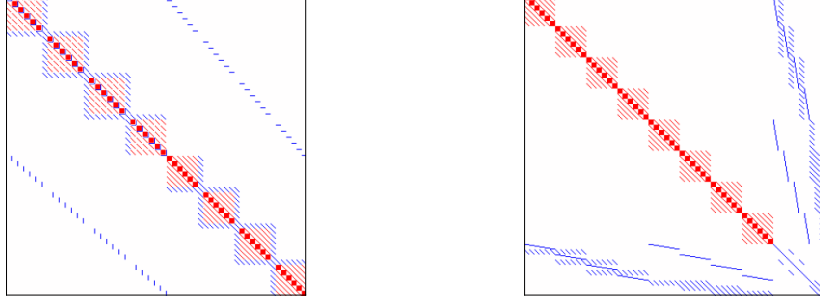
The matrix \mathbf{A} is block-sparse with distinctive structures that are illustrated in Figure 2. It contains relatively dense blocks corresponding to the interactions within subdomain interiors (red in Fig. 2a) bordered by smaller submatrices representing the Neumann derivative approximations between “left” and “right” adjacent subdomains (blue in Fig. 2a). The interactions between “up” and “down” adjacent subdomains are represented by the blue submatrices far from the main diagonal.

Remark 1 (Corner and edge nodes). *The formula (4) does not specify how to handle nodes that lie on corners and edges. Conveniently, it turns out that these nodes can be simply dropped from consideration in any situation where the differential operator in (1) does not involve cross terms of the form $\partial_i \partial_j u$ for $i \neq j$. In situations where cross terms do appear, one reinterpolates the $p \times p$ Chebyshev grid on each face to a $(p-1) \times (p-1)$ grid of Legendre nodes. This way the global system again will involve no discretization nodes on edges or corners. For details, see [9, 14].*

2.2. Static condensation: Elimination of interior nodes

For computational efficiency, we in practice never explicitly form the large sparse matrix \mathbf{A} . Instead we perform a pre-computation on each box that eliminates all nodes that are interior to the box, to form a sparse matrix \mathbf{T} that encodes only interactions between boundary nodes illustrated on a

Matrix Sparsity Patterns



(a) Standard blockwise ordering.

(b) Permuted so that interior nodes go first.

Figure 2: Sparsity pattern of \mathbf{A} for a 2D domain of 4×2 subdomains as in Figure 1. Entries associated with interior nodes are red, and entries associated with boundary nodes are blue. Section 2.2 provides details.

domain in Figure 4. The matrix \mathbf{T} is much smaller than \mathbf{A} (roughly by a factor of p), but is less sparse.

The first step is to partition the discretization nodes into two sets: one that holds all nodes that are interior to a box, and one that holds all nodes that sit on domain boundaries. Let I_i and I_b be two index vectors that identify the two sets. (In other words, I_i marks all diagonal entries that are red in Figure 2, while I_b marks all blue diagonal entries.) Then we partition the linear system (3) accordingly, to obtain

$$\begin{bmatrix} \mathbf{A}_{ii} & \mathbf{A}_{ib} \\ \mathbf{A}_{bi} & \mathbf{A}_{bb} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}_b \end{bmatrix} = \begin{bmatrix} \mathbf{f}_i \\ \mathbf{f}_b \end{bmatrix}, \quad (5)$$

where

$$\mathbf{A}_{ii} = \mathbf{A}(I_i, I_i), \quad \mathbf{A}_{ib} = \mathbf{A}(I_i, I_b), \quad \mathbf{A}_{bi} = \mathbf{A}(I_b, I_i), \quad \mathbf{A}_{bb} = \mathbf{A}(I_b, I_b),$$

and where

$$\mathbf{u}_i = \mathbf{u}(I_i), \quad \mathbf{u}_b = \mathbf{u}(I_b), \quad \mathbf{f}_i = \mathbf{f}(I_i), \quad \mathbf{f}_b = \mathbf{f}(I_b).$$

A key observation here is that the matrix \mathbf{A}_{ii} is block diagonal, since the interior nodes of any single box do not directly communicate with the interior nodes of any other box. This means that we can apply \mathbf{A}_{ii}^{-1} to vectors via

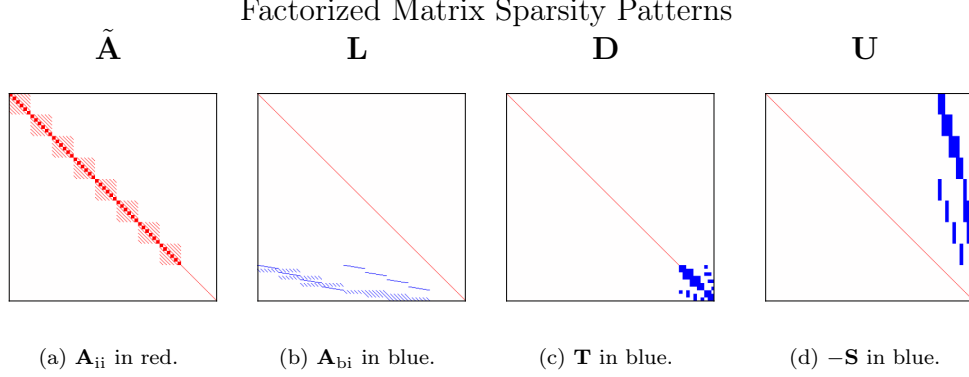


Figure 3: Sparsity patterns of (with $\tilde{\mathbf{A}}$) \mathbf{LDU} factorization for permuted \mathbf{A} . This is for a 2D problem with 4×2 subdomains as in Figure 2. Blocks corresponding to particular submatrices in the factorization are highlighted red or blue.

local computations on each box that are completely unconnected. Another key observation is that even for variable coefficient differential operators, the matrices \mathbf{A}_{bi} and \mathbf{A}_{bb} both consist of repeated, identical submatrices, at least in the case where all subdomains have uniform p . This is because the differential operators they numerically represent - specifically the Neumann derivative to ensure Neumann continuity - are universal across all boxes in the domain. We exploit this invariance to accelerate many of the box computations. A sparsity plot of \mathbf{A} partitioned as in (5) is shown in Figure 2b.

The solver we use relies on the following factorization of the coefficient matrix in (5):

$$\mathbf{A} = \underbrace{\begin{bmatrix} \mathbf{A}_{ii} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}}_{=\tilde{\mathbf{A}}} \underbrace{\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{A}_{bi} & \mathbf{I} \end{bmatrix}}_{=\mathbf{L}} \underbrace{\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_{bb} - \mathbf{A}_{bi}\mathbf{A}_{ii}^{-1}\mathbf{A}_{ib} \end{bmatrix}}_{=\mathbf{D}} \underbrace{\begin{bmatrix} \mathbf{I} & \mathbf{A}_{ii}^{-1}\mathbf{A}_{ib} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}}_{=\mathbf{U}} \quad (6)$$

That (6) holds is easily established by simply multiplying the factors together. It will be convenient to introduce two new matrices

$$\mathbf{S} := -\mathbf{A}_{ii}^{-1}\mathbf{A}_{ib} \quad \text{and} \quad \mathbf{T} := \mathbf{A}_{bb} - \mathbf{A}_{bi}\mathbf{A}_{ii}^{-1}\mathbf{A}_{ib} = \mathbf{A}_{bb} + \mathbf{A}_{bi}\mathbf{S}.$$

From (6), we then easily obtain the inversion formula

$$\mathbf{A}^{-1} = \underbrace{\begin{bmatrix} \mathbf{I} & \mathbf{S} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}}_{=\mathbf{U}^{-1}} \underbrace{\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}^{-1} \end{bmatrix}}_{=\mathbf{D}^{-1}} \underbrace{\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{A}_{bi} & \mathbf{I} \end{bmatrix}}_{=\mathbf{L}^{-1}} \underbrace{\begin{bmatrix} \mathbf{A}_{ii}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}}_{=\tilde{\mathbf{A}}^{-1}}. \quad (7)$$

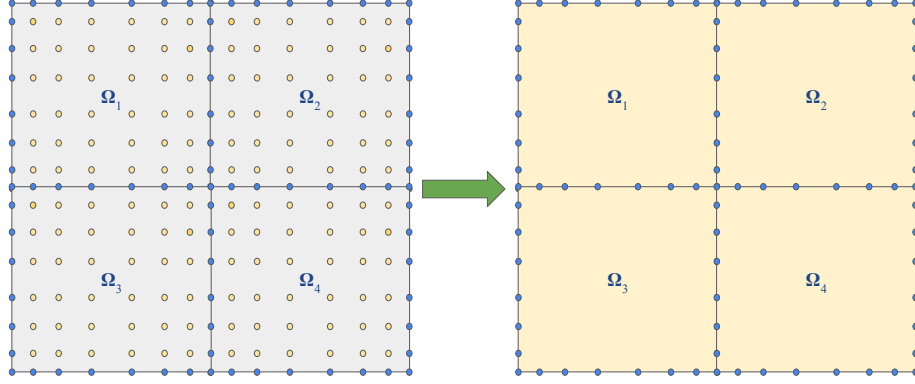


Figure 4: We reduce the problem from one on the entire discretized domain to one on subdomain boundaries only.

Observe that the matrix \mathbf{T} is sparse and is much smaller than the original matrix \mathbf{A} . However, it has a significantly larger ratio of non-zero elements per row than \mathbf{A} . The submatrices \mathbf{S} , \mathbf{A}_{bi} , and \mathbf{A}_{ii}^{-1} are all block matrices, with \mathbf{A}_{bi} in particular consisting of the same repeated block.

To compute the solution \mathbf{u} to the linear system (3) we now need to execute four steps, each one corresponding to one factor in (7):

1. Form the vector $\mathbf{v}_i = \mathbf{A}_{ii}^{-1}\mathbf{u}_i$. Exploit that \mathbf{A}_{ii} is block diagonal to do this through independent local computations, executed via batched linear algebra.
2. Form the vector $\mathbf{g}_b = \mathbf{f}_b - \mathbf{A}_{bi}\mathbf{v}_i$. We observe that \mathbf{A}_{bi} is very sparse and specifically consists of one repeated block, so the vector \mathbf{g}_b can again be formed via a repeated local computation on each box that is accelerated via GPUs (this is a “matrix free” approach).
3. Assemble the block sparse matrix \mathbf{T} through a set of local computations on all boxes. Then solve the global linear system $\mathbf{T}\mathbf{u}_b = \mathbf{g}_b$ using a highly optimized sparse direct solver.
4. Once \mathbf{u}_b is available, we can reconstruct the full solution via the formula $\mathbf{u}_i = \mathbf{v}_i + \mathbf{S}\mathbf{u}_b$.

In situations where the system (1) needs to be solved for multiple right-hand sides, further acceleration can be obtained by pre-computing and storing an LU factorization of the matrix \mathbf{T} . Details of how the algorithm was implemented efficiently are provided in Section 3.

2.3. Discussion

In this manuscript we described the solver from a linear algebraic point of view, as this helps clarify how we attain acceleration via batched linear algebra. However, all steps we take have direct analogues to steps in the original HPS method, as described in [18, Chapter 25]. To be precise, the blocks of \mathbf{T} contain the Dirichlet-to-Neumann operators that are formed in HPS. The matrix \mathbf{S} holds the “solution operators” that reconstruct the solution to (1) at the interior nodes, once the solution at the boundaries has been identified in the global sparse solve.

In the original HPS method the sparse matrix \mathbf{T} is never formed explicitly. Instead, the same solution process is continued hierarchically: first a set of coarser boxes is created by taking the lowest-level boxes and merging them into pairs. This creates a new partition between “interior” and “boundary” boxes which is used to further factorize the matrix \mathbf{T} itself to form an even smaller block sparse matrix \mathbf{T}' , and so on. The solution maps that make up these levels of matrices (\mathbf{T} , \mathbf{T}' , and ongoing) can be formed in parallel within each level.

Highly-optimized multifrontal solvers like MUMPS can use pivoting techniques such as partial threshold pivoting, relaxed pivoting [22], and tournament pivoting [23] to improve numerical stability of a matrix factorization like sparse LU [24]. In theory such pivoting techniques could be utilized on matrices present in previous approaches to HPS to improve their stability. However, encoding solutions for all shared faces into one sparse matrix operator \mathbf{T} allows such pivoting techniques to be applied holistically to a solve for \mathbf{u} on all such points. This means a wider range of pivots can be applied, improving the factorization’s stability while preserving an efficient sparsity pattern. In effect we “look at” a larger portion of the problem at once, allowing us (through the use of a multilevel solver) to perform more optimizations in the factorization process. This approach also enables us to use highly optimized standard libraries in the construction and factorization of \mathbf{T} .

In addition, the non-hierarchical nature of our assembly of \mathbf{T} – relying only on the lowest-level DtN maps and not higher level DtN or solution operators – gives this approach uniformity across the computation of its blocks which is helpful for parallelization. Individually the DtN maps that make up \mathbf{T} are relatively dense matrices, but their entries can be computed on GPUs if available via batched linear algebra routines [25]. This allows us to use GPU accelerators on a significant portion of our method’s computations.

3. Algorithm Description

In this section we elaborate on the implementation of our HPS solver, highlighted by the four steps presented in Section 2.2. We discuss the numerical challenges present in both the local box operations and the formulation of a system for the reduced problem on domain boundaries.

3.1. Local Computations

As shown in Section 2.2, our HPS solver consists of four steps corresponding to the factors $\tilde{\mathbf{A}}^{-1}$, \mathbf{L}^{-1} , \mathbf{D}^{-1} , and \mathbf{U}^{-1} . We can avoid forming these matrices explicitly, and instead apply them in our solvers through matrix-free approaches. The steps to accomplish this fit into two categories of computations.

The first set are batched linear algebra routines that represent $\tilde{\mathbf{A}}$, \mathbf{L} , and \mathbf{U} : the factors are primarily made of submatrices that are block-sparse in structure, such as \mathbf{A}_{ii}^{-1} , \mathbf{S} , and \mathbf{A}_{bi} . Each of these matrices have a limited number of nonzero blocks per row (with each block-row corresponding to solving for \mathbf{u} on one box Ω^τ) that are only in specific indices relative to the row index. Thus we stack these nonzero blocks into three-dimensional arrays of size (number of blocks \times size of each block), with their indices within the matrix inferred from the order in which they are stored. We then apply them to our vector data in a matrix-free approach (using “tensor” objects and operations in PyTorch). Practically speaking, we are representing the action of the matrix-vector product with a series of smaller matrix-vector products corresponding to each nonzero block. By doing this we can use batched linear algebra routines (also present in PyTorch) for the matrix-vector multiplications, which reduce the memory cost and computational time.

Further optimizations can be applied to specific submatrices. For instance, since \mathbf{A}_{ii} is block-diagonal, we can apply \mathbf{A}_{ii}^{-1} to a vector by simply performing a number of independent small system solves. Meanwhile, since \mathbf{A}_{bi} represents part of a Neumann operator used to enforce Neumann continuity on boundary nodes, it is unaffected by our differential operator A so its blocks \mathbf{A}_{bi}^τ are identical across all boxes Ω^τ of our discretization. This means \mathbf{A}_{bi} does not require a stacked 3-dimensional array to represent its multiplication; we may represent it as a single small block applied to each of the subvectors of \mathbf{u} representing a box. This reduces the memory footprint even further than computations with stacked 3D arrays.

The exception to this batched approach is the sparse matrix \mathbf{T} that solves for \mathbf{u} on boundary nodes. While \mathbf{T} is also block-sparse, it must be inverted (or, more accurately, factorized for use in a solver) to apply \mathbf{D}^{-1} , and unlike \mathbf{A}_{ii} it is not block-diagonal. Thus there is no similarly convenient means to invert its components independently, and we instead interface with MUMPS to produce a factorization of it for use in our method.

However, while \mathbf{T}^{-1} cannot easily be applied with a batched matrix-free method in a comparable way to the other submatrices, we may assemble \mathbf{T} in batch before factorizing it. We can write \mathbf{T} as

$$\mathbf{T} = \mathbf{A}_{bb} + \mathbf{A}_{bi}\mathbf{S} = \begin{bmatrix} \mathbf{A}_{bb} & \mathbf{A}_{bi} \end{bmatrix} \begin{bmatrix} \mathbf{I} \\ \mathbf{S} \end{bmatrix}. \quad (8)$$

Since \mathbf{S} itself consists of blocks corresponding to each box Ω^τ that can be neatly represented by a 3D array, we may assemble and store the blocks of \mathbf{T} in a similar fashion, and then apply the blocks of $[\mathbf{A}_{bb}|\mathbf{A}_{bi}]$ to this array (like \mathbf{A}_{bi} , \mathbf{A}_{bb} is identical across blocks so we can represent it as a single small matrix). This gives us the blocks \mathbf{T}^τ that make up \mathbf{T} . From here we can then construct row and column index arrays that align each \mathbf{T}^τ to the subsets of \mathbf{u}_b that reside on $\partial\Omega^\tau$, then construct \mathbf{T} in a sparse format and factorize it with MUMPS, in our case accessed via `petsc4py`.

Remark 2 (Relation to previous HPS implementations). *The submatrices present in our $\tilde{\text{ALDU}}$ factorization of \mathbf{A} each represent operators present in previous formulations of HPS. In particular, \mathbf{A}_{bi} and \mathbf{A}_{bb} are differential matrices approximating the Neumann derivative, the blocks of \mathbf{S} are the solution operators for solving box interiors, and the blocks of \mathbf{T} are the DtN maps consisting of composing Neumann differentiation on solution operators to enforce Neumann continuity. Readers previously familiar with HPS likely notice (8) resembles the formula for constructing a DtN map.*

3.2. Build and Solve stages

Similar to previous HPS implementations, our method can be divided into a build stage (assembling and factorizing the operator \mathbf{T}) and a solve stage (using the LU factors of \mathbf{T} , along with other operators, to solve for \mathbf{u}).

An outline of the build stage is shown in Algorithm 1. In short, we construct the DtN map for each subdomain τ in our partitioned domain, then assemble \mathbf{T} out of these DtNs. Each DtN is of the same size ($\approx 4p \times 4p$ in 2D) and is computed using the same formulation of differentiation operators.

BUILD STAGE

1. Initialize a sparse matrix $\mathbf{T} = \mathbf{0}$.
2. *For each box τ , compute the DtN matrix \mathbf{T}^τ as described in Section 3.1, and add it to the relevant slot of \mathbf{T} .*

for (all boxes τ)
 - Form the local equilibrium matrix \mathbf{A}^τ .
 - Form the local solution operator $\mathbf{S}^\tau = -(\mathbf{A}_{ii}^\tau)^{-1} \mathbf{A}_{ib}^\tau$
 - Form the local DtN matrix $\mathbf{T}^\tau = \mathbf{A}_{bi}^\tau \mathbf{S}^\tau$.
 - Add the matrix \mathbf{T}^τ to the relevant block of \mathbf{T} .

In parallel.
- end for**
3. Form the **LU** factorization of \mathbf{T} .

Algorithm 1: constructs and factorizes \mathbf{T} .

Thus they can be stored in “stacked” three-dimensional arrays and effectively assembled using batched linear algebra routines such as those in PyTorch.

However, a challenge with utilizing GPUs for DtN assembly is memory storage. In 3D each DtN map is size $\approx 6p^2 \times 6p^2$. Building them requires solving matrix systems of the spectral differentiation operators on interior points, which are size $\approx (p-2)^3 \times (p-2)^3$. These matrices grow quickly as p increases (at $p = 20$ they are ≈ 5.7 million and 34 million entries, respectively). For larger problem sizes it is impossible to store all DtN maps on a GPU concurrently. The interior differentiation submatrices in particular do not have to be stored beyond formulating their box’s DtN map, but they are much larger for high p and present considerable memory overhead. To handle this storage requirement we have implemented a two-level scheduler that stores a limited number of DtNs on the GPU at once, and assembles a smaller number of DtNs in batch at a time. For $p \geq 16$ on a V100 GPU the best results occur when only one DtN is assembled at a time. However, we may still store multiple DtNs on the GPU before transferring them collectively to the CPU.

Once the DtNs $\{\mathbf{T}^\tau\}$ are assembled, we add them to the relevant blocks of \mathbf{T} , which is represented as a sparse matrix on the CPU. The row and column indexing arrays used to determine these blocks are also constructed in parallel with PyTorch. Then we factorize \mathbf{T} using MUMPS. This factorization can be

reused for multiple problems with the same differential operator A , avoiding the need to rebuild it.

SOLVE STAGE

1. *Initialize the solution vector: $\mathbf{u} = \mathbf{0}$.*
2. *In a loop over all boxes, form the local solutions \mathbf{v}^τ and the local equivalent load vectors \mathbf{g}^τ :*

for (all boxes τ)

 - Form the matrices \mathbf{A}_{ii}^τ and \mathbf{A}_{bi}^τ
 - Form $\mathbf{u}(I_i^\tau) = (\mathbf{A}_{ii}^\tau)^{-1}\mathbf{f}(I_i^\tau)$. *(So $\mathbf{u}(I_i^\tau)$ holds \mathbf{v}^τ .)*
 - Form $\mathbf{f}(I_b^\tau) = \mathbf{f}(I_b^\tau) - \mathbf{A}_{bi}^\tau\mathbf{u}(I_i^\tau)$. *(So $\mathbf{f}(I_b^\tau)$ holds \mathbf{g}^τ .)*

In parallel.

end for
3. *Solve the global system using a sparse direct solver with the precomputed \mathbf{LU} factorization of \mathbf{T} : $\mathbf{u}(I_b) = \mathbf{U}^{-1}\mathbf{L}^{-1}\mathbf{f}(I_b)$*
4. *Fill in the solution at all the interior nodes:*

for (all boxes τ)

 - $\mathbf{u}(I_i^\tau) = \mathbf{u}(I_i^\tau) + \mathbf{S}^\tau\mathbf{u}(I_b^\tau)$.

In parallel.

end for

Algorithm 2: solve for the PDE given a factorized system \mathbf{T} .

Algorithm 2 outlines the solve stage. Here we compute \mathbf{u} given a differential operator A , resulting boundary system \mathbf{T} , Dirichlet data, and a body load \mathbf{f} . We must reduce \mathbf{f} to the subdomain boundaries before applying it to solve for \mathbf{u} on box boundary points as described in Section 2.2. In addition, after solving for boundary \mathbf{u} we must apply local solution operators \mathbf{S}^τ to \mathbf{u} on each subdomain boundary to obtain \mathbf{u} on the subdomain interiors. Each of these operations can be done in parallel using blocks of \mathbf{A}_{bi} , \mathbf{A}_{ii}^{-1} , and \mathbf{S} . Although these matrices are needed to assemble the local DtN maps, we do not store them after the build stage. For large problems, GPU memory is insufficient, and CPU-GPU transfers are costly, thus it is more efficient to recompute the required block operators during each solve.

The uniformity of DtN assembly and interior subdomain solves is clear when one choice of p is used across all subdomains. However, we could instead select p from a small fixed set of values (perhaps ≈ 6 different choices) to resemble adaptive methods and still retain much of the benefit. In this

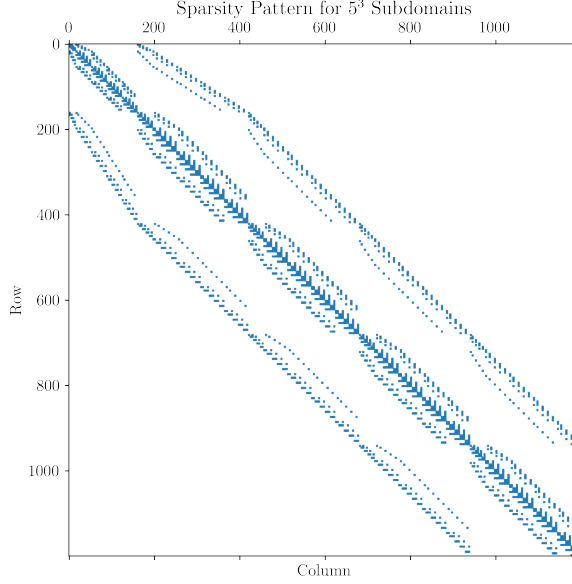


Figure 5: The sparsity pattern of a problem discretized with $5 \times 5 \times 5$ boxes. Some faces have up to 10 faces on the adjacent boxes that are not on domain boundary $\partial\Omega$, but this is independent of the total number of discretization points N .

case we would have separate arrays storing the DtNs for each choice of p , but as long as enough subdomains use each p the corresponding submatrix multiplications can still be batched. The same infrastructure that enables Chebyshev-Legendre interpolation (see Remark 1) can also enable interpolation between adjacent boxes with different p in the assembly of \mathbf{T} , which will be finalized in future work.

3.3. Asymptotic Cost in 3D

The build stage has an asymptotic cost of

$$O(N^2) \text{ (sparse factorization)} + O(p^6 N) \text{ (forming blocks of } \mathbf{T}\text{)}.$$

Based on our numerical experiments the factorization of \mathbf{T} is the most time-consuming portion of our method. Since we are using a multilevel solver for a 3D PDE, it has an asymptotic cost of $O(N^2)$ [18, Chapter 25]. However, this cost is unaffected by the polynomial order p of our discretization. Other

numerical methods such as finite differences with large stencils have a significant prefactor cost in p that can greatly increase the scaling constant in front of this estimate [18, Chapter 20]. Additionally, since \mathbf{T} is relatively small (it has $< \frac{6N}{p} \times \frac{6N}{p}$ total entries with $\approx 11p^2$ nonzero entries per row, with an example shown in Figure 5) it can be stored in memory and reused for multiple solves.

By comparison, increasing p has a substantial effect on the runtime of the DtN computations needed for blocks of \mathbf{T} . This is because each DtN $\mathbf{T}^{(i)}$ requires the solution of a system using a large block of the spectral differentiation matrix $\mathbf{A}^{(i)}$ that is size $\approx p^3 \times p^3$. Nevertheless, we can still assemble DtNs in batch for $p \leq 15$. Even for larger p where batched operations are impractical, building the DtNs in serial is relatively fast as shown in our numerical results. Since DtN assembly is linear in N and can be accelerated with GPUs, a steep cost in p is more tolerable for it than for the factorization of \mathbf{T} .

Our asymptotic cost for the solve stage is

$$O(N^{4/3}) \text{ (sparse solve)} + O(p^6 N) \text{ (box solves with } \mathbf{S}\text{)}.$$

If we stored the solution operators in memory then the box solves are only $O(p^2 N)$, but in practice it is preferred to reformulate the solution operators and avoid having to store and potentially transfer additional memory.

4. Numerical Experiments

In this section we present a series of numerical tests for our HPS solver on various 3D problems, investigating its accuracy and computational performance. Numerical experiments were run on two Intel(R) Xeon(R) Gold 6254 CPUs at 3.10GHz with 768GB RAM, and an Nvidia V100 GPU with 32GB of memory. Unless otherwise noted, the reported errors are the relative ℓ_2 -norm between the computed solution and the known solution as $\|\mathbf{u} - \mathbf{u}_{\text{true}}\|_2 / \|\mathbf{u}_{\text{true}}\|_2$.

4.1. Accuracy results - Poisson and Helmholtz Equations

We investigate the accuracy of our HPS solver on the homogeneous Poisson equation,

$$\begin{aligned} \Delta u &= 0 \text{ on } \Omega = [-1.1, 0.1] \times [1, 2] \times [1.2, 2.2] \\ u(\mathbf{x}) &= \frac{1}{4\pi\|\mathbf{x}\|_2} \text{ on } \partial\Omega. \end{aligned} \tag{9}$$

Relative Errors for the Poisson and Helmholtz (10 ppw) Equations

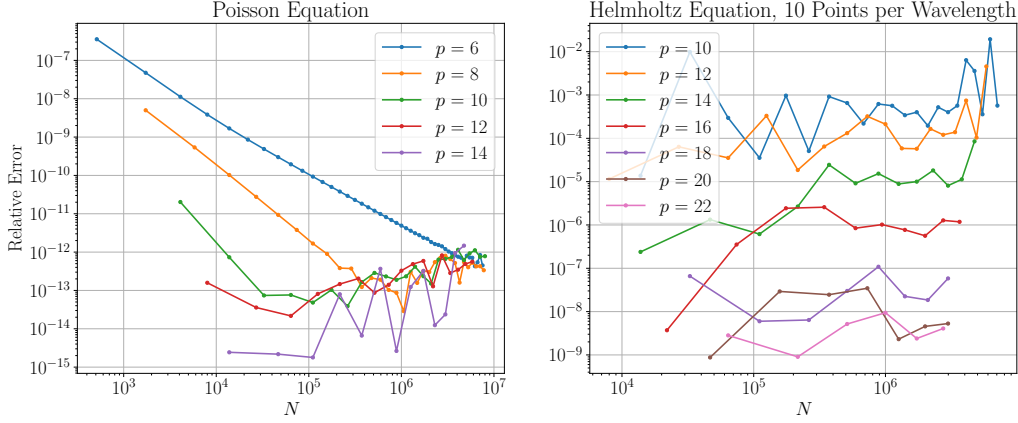


Figure 6: (*left*) Accuracy of our solver for the Poisson Equation. (*right*) Accuracy for the Helmholtz equation fixed to 10 points-per-wavelength. Both cases feature h -refinement with each line marking a fixed p .

Here the Dirichlet boundary condition is defined so that it has a solution in its Green's function. We compute numerical solutions to this equation across a range of block sizes h and polynomial orders p , comparing their relative ℓ^2 -norm error to the Green's function. The accuracy results are in Figure 6. We notice convergence rates of $\approx h^7$ for $p = 6$ and $\approx h^{12}$ for $p = 8$ until the relative error reaches $\approx 10^{-13}$, at which point it flattens (for larger p the error reaches 10^{-13} before a convergence rate can be fairly estimated). Subsequent refinement does not consistently improve accuracy, but it does not lead to an increase in error beyond $\approx 10^{-12}$ either. This suggests our solver is relatively stable.

To assess our method's accuracy on oscillatory problems we also evaluate it on the homogeneous Helmholtz equation with its Green's function used as a manufactured solution with wavenumber κ ,

$$\begin{aligned} \Delta u + \kappa^2 u &= 0 \text{ on } \Omega = [-1.1, 0.1] \times [1, 2] \times [1.2, 2.2] \\ u(\mathbf{x}) &= \frac{\cos(\kappa|\mathbf{x}|)}{4\pi|\mathbf{x}|} \text{ on } \partial\Omega. \end{aligned} \quad (10)$$

In Figure (7) we show relative ℓ^2 -norm errors for our solver on Equation (10) with fixed wavenumbers $\kappa = 16$ and $\kappa = 30$, which correspond to roughly 2.5 and 4.5 wavelengths on the domain, respectively. We see comparable convergence results to the Poisson equation, though the lower bound on relative

Relative Errors for Helmholtz Equation with Fixed κ

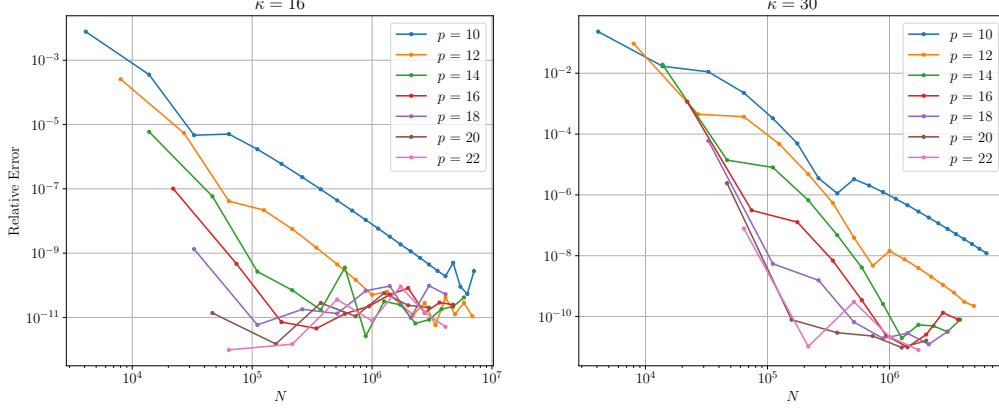


Figure 7: Accuracy of our solver for the Helmholtz equation set to wavenumbers $\kappa = 16$ (≈ 2.5 wavelengths on the domain) and $\kappa = 30$ (≈ 4.5 wavelengths).

error is somewhat higher for each problem ($\approx 10^{-11}$ for $\kappa = 16$ and $\approx 10^{-10}$ for $\kappa = 30$). However, increasing the resolution further does not cause the error to diverge. Along with the results for Equation (9), this suggests our method remains stable when problems are over-resolved. Figure 6 shows the relative error for the Helmholtz equation with a varying κ , set to provide ≈ 10 discretization points per wavelength. We see higher p corresponds to a better numerical accuracy, while increasing h does not improve results (since the wavenumber is also increased), but errors remain relatively stable. Since our method maintains accuracy while increasing the total degrees of freedom linearly with κ^3 , it does not appear to suffer from the pollution effect.

We then consider the gravity Helmholtz equation in an example similar to one in [26],

$$\begin{aligned} \Delta u + \kappa^2(1 - x_3)u &= -1 \text{ on } \Omega = [1.1, 2.1] \times [-1, 0] \times [-1.2, -0.2] \\ u(\mathbf{x}) &= 0 \text{ on } \partial\Omega. \end{aligned} \quad (11)$$

This equation adds a spatially-varying component to the wavenumber variation based on x_3 that is analogous to the effect of a gravitational field, hence the name [27]. Since $1 < 1 - x_3 < 2.2$, the variation is greater through Ω than in the constant Helmholtz equation with the same κ . Unlike the previous Poisson and Helmholtz equations, equation (11) does not have a manufactured solution. It also produces relatively sharp gradients near the domain boundaries since there is a zero Dirichlet boundary condition, and it has a

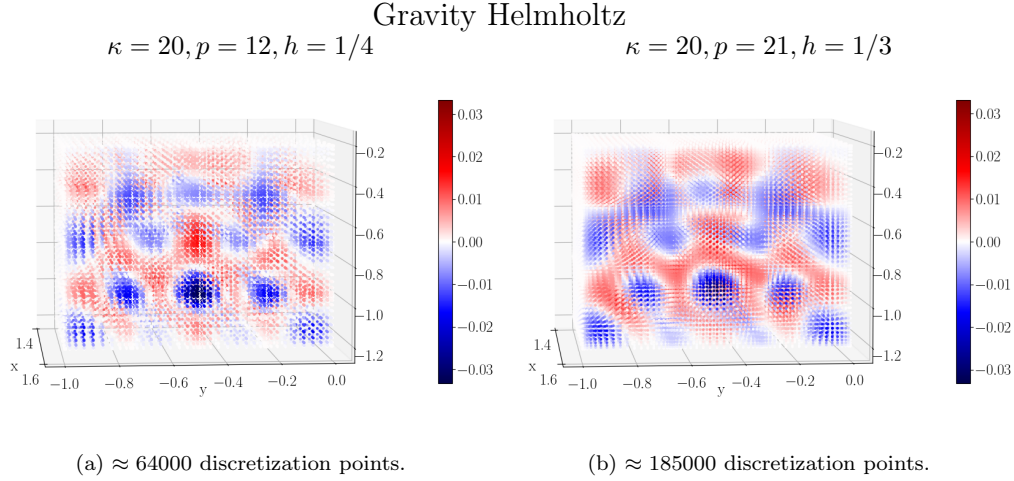


Figure 8: Plots of solutions to the gravity Helmholtz equation as described in Equation (11). A cross-section has been taken through the x -axis. We see consistent results across different p and h .

variable coefficient in spatial coordinates (specifically x_3). We investigate our solver’s accuracy by computing the relative ℓ^2 -norm errors of our results to an over-resolved solution at select points. Figure 8 shows the solution of this problem using select p and h values, while Figure 9 shows these errors in the case of p -refinement for a range of values of h . We can see our solver shows convergence for each h .

4.2. Speed and scaling

To evaluate the computational complexity of our HPS solver, we measure the execution times of four components: the assembly of the DtN operators used for the build stage, the factorization time for our sparse matrix constructed by the DtNs, the solve time for the factorized system, and the time for batched solves of the box subdomain interiors. We know the asymptotic cost of constructing the DtN operators is $O(p^6 N)$ as shown in Equation (9). Figure (10) illustrates the run times for the build stage in solving the Poisson and Helmholtz equations detailed in Section 4.1, using batched linear algebra routines on an Nvidia V100 GPU. As expected, we observe a linear increase in runtime with h -refinement (i.e. a higher N with fixed p) for the DtN assembly. The constant for the linear scaling in h , though, increases sharply with higher values of p .

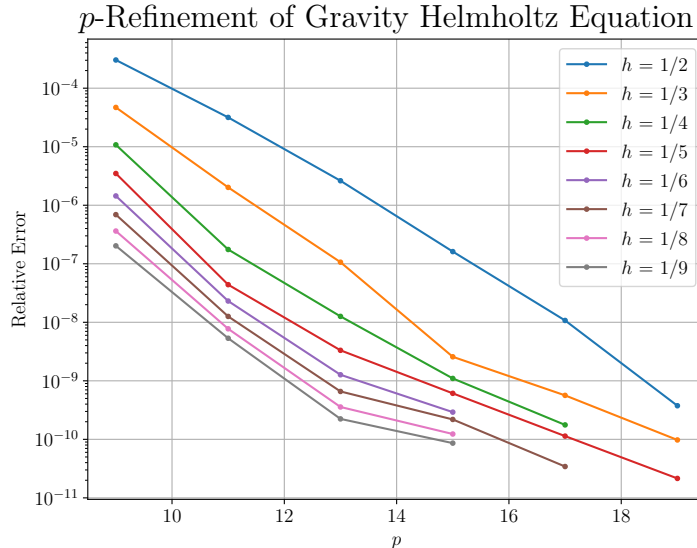


Figure 9: Convergence of our solver on the gravity Helmholtz equation with $\kappa = 12$, relative to a highly-resolved solution. Here h is fixed and resolution is increased through p -refinement.

In Figure 10 we also see the time to factorize the resulting sparse matrix for the build stage, using MUMPS. The expected asymptotic cost is $O(N^2)$. However, our results closely resemble a lower cost of $O(N^{3/2})$ for N up to ≈ 8 million. Overall the factorization time dominates the build stage of our method, given its more costly scaling in N and the reliance on CPUs. Lastly in Figure 11 we see the runtime for the solves of factorized \mathbf{T} for the box boundaries and the batched box interior solves, respectively. Although the box boundary solve has a higher asymptotic cost in N , in practice it has a shorter runtime, especially for large p . The box interior solves follow a linear scaling similar to the batched DtN construction.

4.3. Curved and non-rectangular domains

The HPS method has been extended to other problems such as surface PDEs [28], inverse scattering problems [29], and more generally non-rectangular domains. We can apply our HPS solver to non-rectangular domain geometries through the use of an analytic parameterization between the domain we wish to model, such as a sinusoidal curve along the y -axis shown in Figure 12a, and a rectangular reference domain. These parameter maps extend the versatility of the solver, but feature multiple challenges in their

Build Stage Times

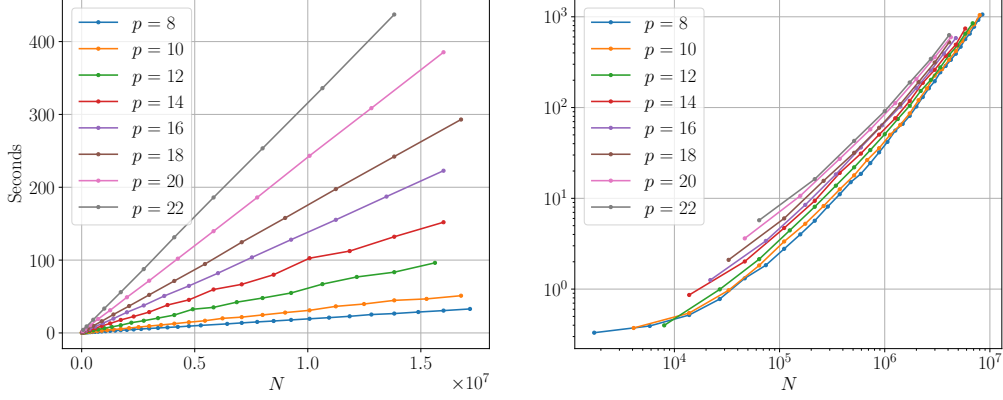


Figure 10: Time to formulate the DtN maps necessary for the sparse system \mathbf{T} (*left*) and factorize \mathbf{T} (*right*).

implementation: they require variable coefficients for our differential operators; they may use mixed second order differential terms, i.e. $\partial^2 u / (\partial x_i \partial x_j)$ where $i \neq j$; and (depending on the domain) may have areas that are near-singular, such as around the sharp point $(x, y) \approx (2, 0)$ in Figure 12a. We investigate the accuracy of our solver for the Helmholtz equation as shown in Equation (10), but now on the sinusoidal domain

$$\Psi = \left\{ x_1, \frac{x_2}{\psi(x_1)}, x_3 \text{ for } (x_1, x_2, x_3) \in \Omega \right\} \text{ where } \psi(z) = 1 - \frac{1}{4} \sin(6z). \quad (12)$$

We map Eq. (10) on Ψ to a cubic reference domain $\Omega = [1.1, 2.1] \times [-1, 0] \times [-1.2, 0.2]$ using the parameter map

$$\begin{aligned} -\frac{\partial^2 u}{\partial x_1^2} - \left(\left(\frac{\psi'(x_1)x_2}{\psi(x_1)} \right)^2 + \psi(x_1)^2 \right) \frac{\partial^2 u}{\partial x_2^2} - \frac{\partial^2 u}{\partial x_3^2} - 2 \frac{\psi'(x_1)x_2}{\psi(x_1)} \frac{\partial^2 u}{\partial x_1 \partial x_2} \\ - \frac{\psi''(x_1)x_2}{\psi(x_1)} \frac{\partial u}{\partial x_2} - \kappa u = 0, \quad (x_1, x_2, x_3) \in \Omega. \end{aligned} \quad (13)$$

We test this problem with a manufactured solution in the Green's function of the Helmholtz equation, with corresponding Dirichlet data matching our test in Section 4.1. Convergence studies for this problem with $\kappa = 16$ and $\kappa = 30$ are shown in Figure 13. We observe a similar convergence pattern to

T Solve Times

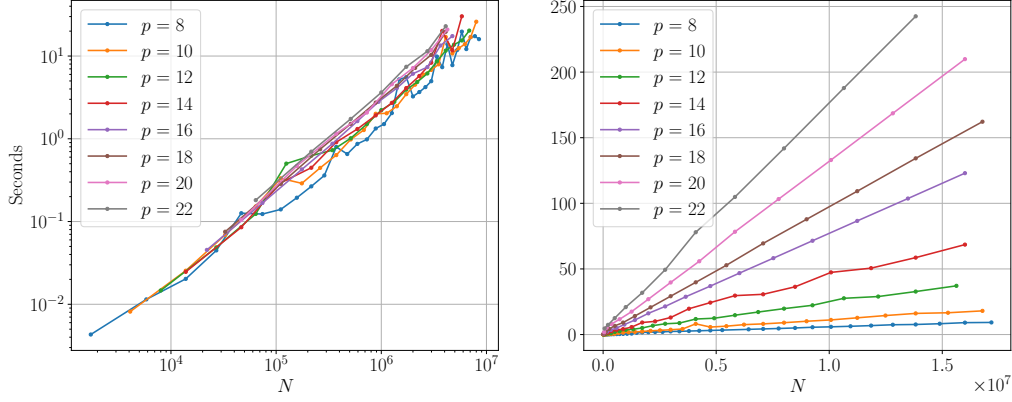


Figure 11: Time to solve a problem using the factorized sparse system \mathbf{T} (*left*) and apply the interior box solves (*right*).

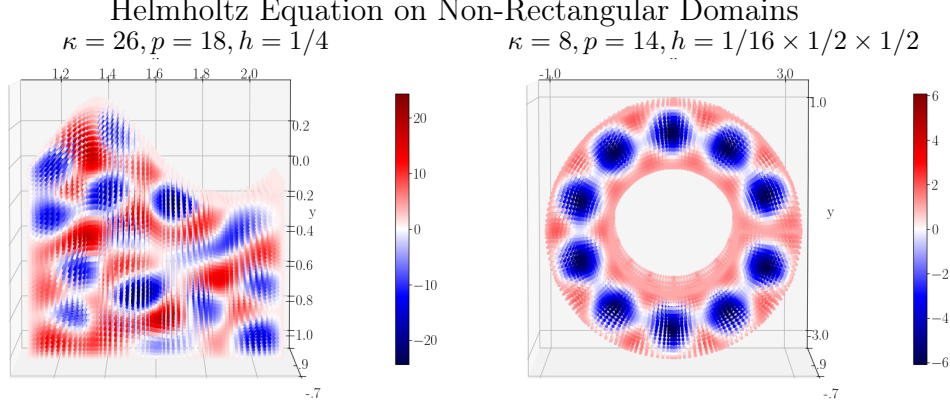
the Helmholtz test on a cubic domain, although the speed of convergence is somewhat slower. A plot of this problem for wavenumber $\kappa = 26$ with a zero Dirichlet boundary condition and body load $f = 1$ is shown in Figure 12a, to provide an example that does not have a manufactured solution.

We also investigate our solver's accuracy on the Helmholtz equation when the domain is a three-dimensional annulus, with an example with zero Dirichlet boundary condition and $f = 1$ shown in Figure 12b. We use an analytic parametrization to map this problem to a rectangular reference domain with higher resolution in the y -axis and a periodic boundary condition. Relative errors in p and N are comparable to those for the problem on the sinusoidal domain.

4.4. Parabolic PDEs with time stepping

Consider a convection-diffusion equation modeling the concentration of a contaminant u in a fluid, over time interval $[0, T = 5]$ on $\Omega = [-0.5, 0.5]^3$, with diffusivity $\kappa = 10^{-4}$ and circular horizontal velocity field $\mathbf{b}(\mathbf{x})$:

$$\begin{aligned} \mathbf{b}(x_1, x_2, x_3) &= (-\cos(x_1) \sin(x_2)x_3, \sin(x_1) \cos(x_2)x_3, 0) \\ \frac{\partial u(\mathbf{x}, t)}{\partial t} &= \kappa \Delta u - \nabla \cdot (\mathbf{b}(\mathbf{x})u) = Au, & (x, t) &\in \Omega \times [0, T], \\ u(\mathbf{x}, t) &= 0, & (x, t) &\in \partial\Omega \times [0, T]. \end{aligned} \tag{14}$$



(a) On a curved domain with a wavy edge, $\approx 2.6 \times 10^5$ discretization points. (b) On an annulus, with $\approx 6.9 \times 10^4$ discretization points.

Figure 12: Helmholtz equation on two non-rectangular domains, both with a Dirichlet boundary condition of 1. A cross section has been taken along the z axis.

Here A is the linear partial differential operator where $Au = (\kappa \Delta u - \nabla \cdot (\mathbf{b}u))$. Equation (14) is a parabolic PDE that requires both temporal and spatial discretization. We will apply a temporal discretization through the Crank-Nicolson method, which gives us

$$\begin{aligned} \frac{u^{n+1} - u^n}{\Delta t} &= \frac{1}{2} A u^{n+1} + \frac{1}{2} A u^n \\ \Rightarrow \left(I - \frac{\Delta t}{2} A \right) u^{n+1} &= \left(I + \frac{\Delta t}{2} A \right) u^n. \end{aligned} \quad (15)$$

The left-hand-side of (15) is itself a linear elliptic partial differential operator, thus we may discretize it using our HPS method and produce both box solution and DtN maps as well as the system \mathbf{T} . Moreover, the same PDO is used for every time step in the computation. Thus the build stage of our method need only be applied once - we can use the same factorized \mathbf{T} for every time step, supplying \mathbf{u}^n from the previous timestep as the body load. However, we must also multiply \mathbf{u}^n by a similar but distinct operator, $I + (\Delta t/2)A$. This is applied to both the box solves (by setting $\mathbf{f} = (\mathbf{I} + (\Delta t/2)\mathbf{A})[I, :] \mathbf{u}^n$ for our interior box solve), and the right-hand-side for the sparse system solve in \mathbf{T} . Previous works such as [20, 30] have explored this concept for using HPS to solve time-dependent problems.

Relative Errors for Helmholtz Equation on a Curved Domain

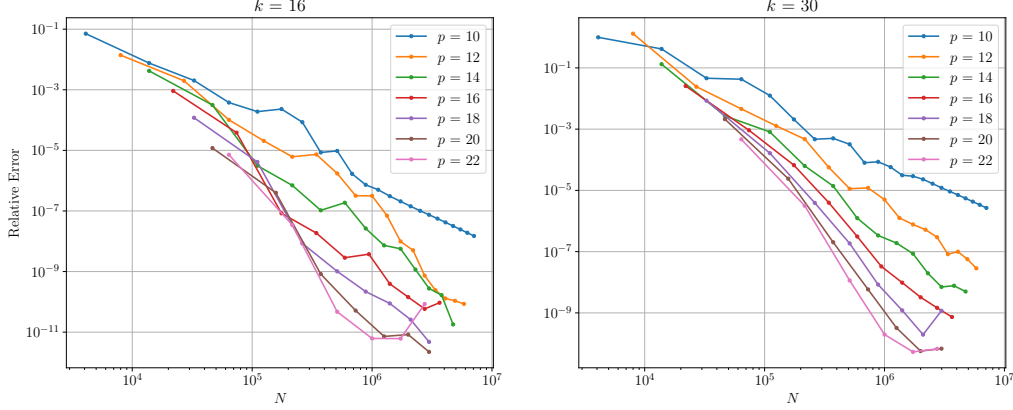


Figure 13: Accuracy of our solver for the Helmholtz equation on a curved domain shown in Fig. 12a set to wavenumbers $\kappa = 16$ (≈ 2.5 wavelengths on domain) and $\kappa = 30$ (≈ 4.5 wavelengths).

We investigate the use of our HPS method in solving the PDE in Equation (14) with the initial condition $u^0 = \exp((x_1^2 + (x_2 + 0.3)^2 + x_3^2)/0.002)$. This creates a contaminant density positioned slightly below the center of Ω along the y axis. Thus the contaminant u will be carried by the velocity field in a counterclockwise pattern. The progress of this flow can be seen in Figure 14. As expected, the contaminant is beginning to move in the pattern determined by $\mathbf{b}(\mathbf{x})$ as the convective term.

For modeling a parabolic PDE we must consider the impact of both spatial and temporal refinement. The Crank-Nicolson method is numerically stable, but it is only a second-order method. Thus refining the size of timestep Δt may have a more significant effect on numerical results than refining p or h in the HPS method. This is explored in Figure 15. Reducing Δt does increase the method's overall runtime, but since the same sparse system \mathbf{T} is used for each step we bypass the build stage of the method, including the costly factorization of \mathbf{T} . In practice HPS would be better combined with a higher order time stepping scheme such as certain Runge-Kutta methods, but as long as the underlying partial differential operator is linear we may leverage the advantage of applying the build stage only once.

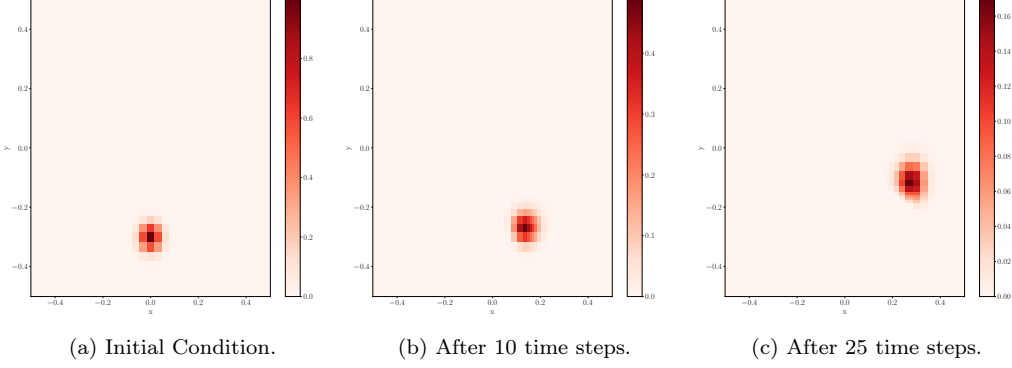


Figure 14: Progression of contaminant \mathbf{u} in a circular flow up to $t = 2.5$. This is a cross-section of a three-dimensional space, so vertical diffusion also affects the result. We see \mathbf{u} follows an expected counterclockwise pattern, continued in Figure 15.

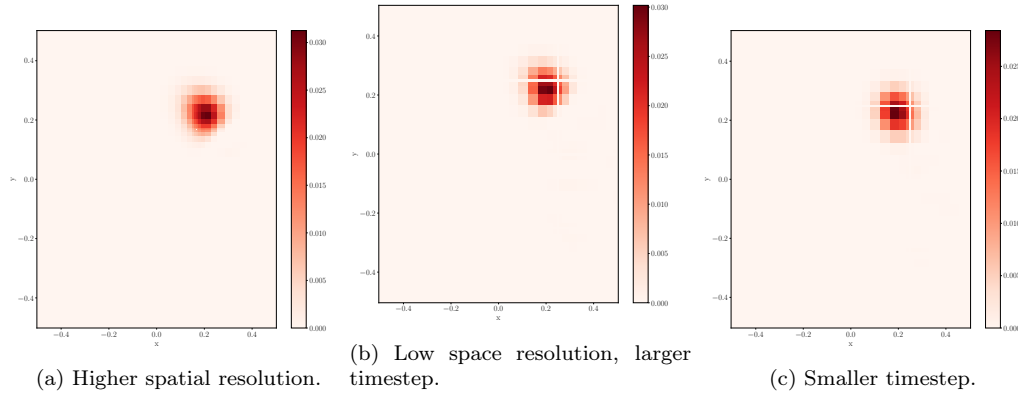


Figure 15: Comparison of \mathbf{u} at $t = 5$ in three regimes: lower spatial resolution and larger time step in the middle, with higher space or time resolutions on either side.

5. Conclusion and Future Work

We present a Hierarchical Poincaré–Steklov scheme tailored for modern hardware, combining GPU-accelerated local solves with a black-box multifrontal sparse solver to deliver both numerical stability and scalability. By casting the global interface system in a form amenable to multifrontal factorization, we decouple dense local computations from the sparse interface solve. In our two-level approach, batch-processed linear algebra on the GPU first factorizes the block-diagonal subdomain systems; then a sparse system on the remaining interfaces is assembled and factorized using black-box sparse solvers, yielding enhanced performance and robustness.

The method presented can readily be adapted to handle other versions of HPS, including those relying on impedance-to-impedance maps [31] and those that use expansion coefficients to represent interface functions [26]. It could in principle also be used in situations where adaptivity in either the local discretization order p or in the mesh size is used [32, 33], or with formulations that exploit internal structure in the local spectral differentiation matrices [34]. However, it may be a delicate matter to maintain the very high efficiency provided by batched linear algebra in our implementation.

Rank-structured compression offers a means of reducing the cost of the global interface solve. Dense Schur complements arising from the coarse level of a sparse factorization are often rank-compressible [18, 35, 36, 37] and can be recovered using black-box randomized methods [38, 39, 40, 41, 42, 43]. Slab-based domain decomposition schemes [44, 45, 46] are particularly compelling because they confine most operations to local subdomains, so that global coupling is restricted to a small, coarse-scale system — simplifying parallelism and enabling effective use of rank-structured compression compared to nested-dissection approaches. Pursuing these ideas will allow the HPS method to scale efficiently to larger problems while preserving robustness for complex physical systems.

6. Code Availability

A code repository of the HPS solver described in this paper is available here: <https://doi.org/10.5281/zenodo.16379968> [47].

7. Conflicts of Interest and Acknowledgments

The authors have no conflicts of interest to report.

This work was supported by the Office of Naval Research (N00014-18-1-2354), by the National Science Foundation (DMS-1952735 and DMS-2313434), and by the Department of Energy ASCR (DE-SC0022251).

References

- [1] Anna Yesypenko and Per-Gunnar Martinsson. GPU optimizations for the hierarchical Poincaré-Steklov scheme. In *International Conference on Domain Decomposition Methods*, pages 519–528. Springer, 2022.
- [2] P.G. Martinsson. A direct solver for variable coefficient elliptic PDEs discretized via a composite spectral collocation method. *Journal of Computational Physics*, 242(0):460 – 479, 2013.
- [3] Patrick R Amestoy, Iain S Duff, Jean-Yves L’Excellent, and Jacko Koster. MUMPS: a general purpose distributed memory sparse solver. In *International Workshop on Applied Parallel Computing*, pages 121–130. Springer, 2000.
- [4] Quang Dinh and Yves Marechal. Toward real-time finite-element simulation on GPU. *IEEE Transactions on Magnetics*, 52(3):1–4, 2015.
- [5] Serban Georgescu, Peter Chow, and Hiroshi Okuda. GPU acceleration for FEM-based structural analysis. *Archives of Computational Methods in Engineering*, 20:111–121, 2013.
- [6] Andreas Klöckner, Timothy Warburton, and Jan S Hesthaven. High-order discontinuous Galerkin methods by GPU metaprogramming. In *GPU Solutions to Multi-scale Problems in Science and Engineering*, pages 353–374. Springer, 2013.
- [7] Axel Modave, Amik St-Cyr, Wim A Mulder, and Tim Warburton. A nodal discontinuous Galerkin method for reverse-time migration on GPU clusters. *Geophysical Journal International*, 203(2):1419–1435, 2015.
- [8] Sagar Imambi, Kolla Bhanu Prakash, and GR Kanagachidambaresan. PyTorch. *Programming with TensorFlow: solution for edge computing applications*, pages 87–104, 2021.

- [9] Sijia Hao and Per-Gunnar Martinsson. A direct solver for elliptic PDEs in three dimensions based on hierarchical merging of Poincaré-Steklov operators. *Journal of Computational and Applied Mathematics*, 308:419–434, 2016.
- [10] Adrianna Gillman and Per-Gunnar Martinsson. An $O(N)$ algorithm for constructing the solution operator to 2D elliptic boundary value problems in the absence of body loads. *Advances in Computational Mathematics*, 40:773–796, 2014.
- [11] Natalie N Beams, Adrianna Gillman, and Russell J Hewett. A parallel shared-memory implementation of a high-order accurate solution technique for variable coefficient Helmholtz problems. *Computers & Mathematics with Applications*, 79(4):996–1011, 2020.
- [12] José Pablo Lucero Lorca, Natalie Beams, Damien Beecroft, and Adrianna Gillman. An iterative solver for the HPS discretization applied to three dimensional Helmholtz problems. *SIAM Journal on Scientific Computing*, 46(1):A80–A104, 2024.
- [13] Adrianna Gillman, AlexH. Barnett, and Per-Gunnar Martinsson. A spectrally accurate direct solution technique for frequency-domain scattering problems with variable media. *BIT Numerical Mathematics*, 55(1):141–170, 2015.
- [14] PG Martinsson. The hierarchical Poincaré-Steklov (HPS) solver for elliptic PDEs: A tutorial. *arXiv preprint arXiv:1506.01308*, 2015.
- [15] Ivo Babuška, Frank Ihlenburg, Ellen T Paik, and Stefan A Sauter. A generalized finite element method for solving the Helmholtz equation in two dimensions with minimal pollution. *Computer methods in applied mechanics and engineering*, 128(3-4):325–359, 1995.
- [16] Ivo M Babuska and Stefan A Sauter. Is the pollution effect of the FEM avoidable for the Helmholtz equation considering high wave numbers? *SIAM Journal on numerical analysis*, 34(6):2392–2423, 1997.
- [17] Jeffrey Galkowski and Euan A Spence. Does the Helmholtz boundary element method suffer from the pollution effect? *Siam Review*, 65(3):806–828, 2023.

- [18] Per-Gunnar Martinsson. *Fast direct solvers for elliptic PDEs*. SIAM, 2019.
- [19] Timothy A Davis. *Direct methods for sparse linear systems*. SIAM, 2006.
- [20] Ke Chen, Daniel Appelö, Tracy Babb, and Per-Gunnar Martinsson. Fast and high-order approximation of parabolic equations using hierarchical direct solvers and implicit Runge-Kutta methods. *Communications on Applied Mathematics and Computation*, pages 1–21, 2024.
- [21] L.N. Trefethen. *Spectral Methods in Matlab*. SIAM, Philadelphia, 2000.
- [22] Iain S Duff and Stéphane Pralet. Towards stable mixed pivoting strategies for the sequential and parallel solution of sparse symmetric indefinite systems. *SIAM Journal on Matrix Analysis and Applications*, 29(3):1007–1024, 2007.
- [23] Laura Grigori, James W Demmel, and Hua Xiang. Communication avoiding Gaussian elimination. In *SC’08: Proceedings of the 2008 ACM/IEEE Conference on Supercomputing*, pages 1–12. IEEE, 2008.
- [24] P. Amestoy, J.-Y. L’Excellent, T. Mary, and C. Puglisi. MUMPS: Multifrontal massively parallel solver for the direct solution of sparse linear equations. In *EMS/ECMI Lanczos prize, ECM conference, Sevilla*, pages 15–19, July 2024.
- [25] Jack Dongarra, Sven Hammarling, Nicholas J Higham, Samuel D Relton, and Mawussi Zounon. Optimized batched linear algebra for modern architectures. In *European Conference on Parallel Processing*, pages 511–522. Springer, 2017.
- [26] Daniel Fortunato, Nicholas Hale, and Alex Townsend. The ultraspherical spectral element method. *Journal of Computational Physics*, 436:110087, 2021.
- [27] Alex H Barnett, Bradley J Nelson, and J Matthew Mahoney. High-order boundary integral equation solution of high frequency wave scattering from obstacles in an unbounded linearly stratified medium. *Journal of Computational Physics*, 297:407–426, 2015.

- [28] Daniel Fortunato. A high-order fast direct solver for surface PDEs. *SIAM Journal on Scientific Computing*, 46(4):A2582–A2606, 2024.
- [29] Carlos Borges, Adrianna Gillman, and Leslie Greengard. High resolution inverse scattering in two dimensions using recursive linearization. *SIAM Journal on Imaging Sciences*, 10(2):641–664, 2017.
- [30] Tracy Babb, Per-Gunnar Martinsson, and Daniel Appelö. Hps accelerated spectral solvers for time dependent problems: Part i, algorithms. In *Spectral and High Order Methods for Partial Differential Equations ICOSAHOM 2018: Selected Papers from the ICOSAHOM Conference, London, UK, July 9-13, 2018*, pages 131–141. Springer International Publishing, 2020.
- [31] Adrianna Gillman, Alex H Barnett, and Per-Gunnar Martinsson. A spectrally accurate direct solution technique for frequency-domain scattering problems with variable media. *BIT Numerical Mathematics*, 55:141–170, 2015.
- [32] Damyn Chipman, Donna Calhoun, and Carsten Burstedde. A fast direct solver for elliptic PDEs on a hierarchy of adaptively refined quadrees. *arXiv preprint arXiv:2402.14936*, 2024.
- [33] Peter Geldermans and Adrianna Gillman. An adaptive high order direct solution technique for elliptic boundary value problems. *SIAM Journal on Scientific Computing*, 41(1):A292–A315, 2019.
- [34] J Aurentz. Fast algorithms for spectral differentiation matrices. *Electron. Trans. Numer. Anal.*, 44:281–288, 2015.
- [35] Steffen Börm, Lars Grasedyck, and Wolfgang Hackbusch. Introduction to hierarchical matrices with applications. *Engineering analysis with boundary elements*, 27(5):405–422, 2003.
- [36] Sabine Le Borne, Lars Grasedyck, and Ronald Kriemann. Domain-decomposition based \mathcal{H} -LU preconditioners. In *Domain decomposition methods in science and engineering XVI*, volume 55 of *Lect. Notes Comput. Sci. Eng.*, pages 667–674. Springer, Berlin, 2007.

- [37] Jianlin Xia. Efficient structured multifrontal factorization for general large sparse matrices. *SIAM Journal on Scientific Computing*, 35(2):A832–A860, 2013.
- [38] James Levitt and Per-Gunnar Martinsson. Linear-complexity black-box randomized compression of rank-structured matrices. *SIAM Journal on Scientific Computing*, 46(3):A1747–A1763, 2024.
- [39] L. Lin, J. Lu, and L. Ying. Fast construction of hierarchical matrix representation from matrix-vector multiplication. *Journal of Computational Physics*, 230(10):4071 – 4087, 2011.
- [40] P.G. Martinsson. A fast randomized algorithm for computing a hierarchically semiseparable representation of a matrix. *SIAM Journal on Matrix Analysis and Applications*, 32(4):1251–1274, 2011.
- [41] Katherine J. Pearce, Anna Yesypenko, James Levitt, and Per-Gunnar Martinsson. Randomized rank-structured matrix compression by tagging, 2025.
- [42] Jianlin Xia. Randomized sparse direct solvers. *SIAM Journal on Matrix Analysis and Applications*, 34(1):197–227, 2013.
- [43] Anna Yesypenko and Per-Gunnar Martinsson. Randomized strong recursive skeletonization: Simultaneous compression and LU factorization of hierarchical matrices using matrix-vector products, 2023.
- [44] Björn Engquist and Lexing Ying. Sweeping preconditioner for the Helmholtz equation: hierarchical matrix representation. *Communications on pure and applied mathematics*, 64(5):697–735, 2011.
- [45] Martin J Gander and Hui Zhang. A class of iterative solvers for the Helmholtz equation: Factorizations, sweeping preconditioners, source transfer, single layer potentials, polarized traces, and optimized schwarz methods. *Siam Review*, 61(1):3–76, 2019.
- [46] Anna Yesypenko and Per-Gunnar Martinsson. SlabLU: a two-level sparse direct solver for elliptic PDEs. *Advances in Computational Mathematics*, 50(4):90, 2024.
- [47] Joseph Kump and Anna Yesypenko. hpslib/hpsmultidomain: Initial release, July 2025. 10.5281/zenodo.16379968.